

SHOTGUN ASSEMBLY OF LABELED GRAPHS

Elchanan Mossel *

Nathan Ross †

September 29, 2015

Abstract

We consider the problem of reconstructing graphs or labeled graphs from neighborhoods of a given radius r . Special instances of this problem include DNA shotgun assembly, neural network reconstruction, and assembling random jigsaw puzzles. We provide some necessary and some sufficient conditions for correct recovery both in combinatorial terms and for some generative models including random labelings of lattices, Erdős-Rényi random graphs, and the random jigsaw puzzle model. Many open problems and conjectures are provided.

1 Introduction

In this paper we study the problem of inferring a (labeled) graph from a collection of radius r (labeled) neighborhoods of the graph. In particular we ask how large r must be to ensure that a given randomly generated graph with labels can be uniquely identified (up to isomorphism) by its r -neighborhoods. Note that if the neighborhoods are too small then identifiability may be impossible: if $r = 1$ and all of the vertex labels are the same, then the graph is only identifiable from its 1-neighborhoods if the degree sequence determines a unique graph. As far as we know graph shotgun assembly for generative models has not been considered before in the level of generality considered here. Some motivating examples include:

- DNA shotgun assembly: the goal is to reconstruct a DNA sequence from “shotgunned” stretches of the sequence. The theoretical version of this problem is graph shotgun assembly of a line graph with each vertex corresponding to a site in the genome, and so is labeled with an A, C, G , or T standing for the nucleotides making up DNA. The neighborhoods are strings of adjacent vertices of length r , which are referred to as “reads”. Shotgun assembly is one of the major techniques for reading DNA sequences and so the theoretical problem is already well understood. A main question is to determine how large r has to be to reconstruct the sequence with good probability under different models of vertex labeling, see e.g., [Arratia et al., 1996], [Dyer et al., 1994], and [Motahari et al., 2013] and references therein.
- Reconstructing neural networks: recent work in applied neuroscience identifies graph shotgun assembly as an important problem for reconstructing neural networks; the goal is to reconstruct a big neural network from subnetworks that are observed in experiments [Soudry et al., 2013].
- The *random jigsaw puzzle problem*. Consider a jigsaw puzzle of size $n \times n$ where where the border between every two adjacent pieces is drawn uniformly at random using one in q shapes

*University of California, Berkeley and University of Pennsylvania; mossel@wharton.upenn.edu

†University of Melbourne; nathan.ross@unimelb.edu.au

of interfaces which we call “jigs.” How large should q be so that the puzzle can be recovered uniquely? How can this be done efficiently?

The problem considered here is most closely related to the famous *reconstruction conjecture* in combinatorics [Kelly, 1957] [Harary, 1974] which can be stated as follows: a graph G on at least 3 vertices is uniquely determined by the multi-set of all vertex-deleted subgraphs of G . Here a vertex deleted subgraph of G is a graph induced on all the vertices of G but one. In this paper we are interested in reconstructing (labeled) graphs from seemingly less information: given the graph we assume that we are given all (labeled) radius r neighborhoods in the graph. While the information is more localized, we make the additional assumption that either the graph structure or the labels are random. This makes the problem easier in comparison to the reconstruction conjecture. Indeed we show that for some popular random graphs models reconstruction is possible from relatively small neighborhoods.

The graph shotgun problem is also related to the graph isomorphism problem [Babai et al., 1980]. It is a famous open problem to determine the complexity of graph isomorphism. In fact, one may consider a variant of the graph isomorphism problem in our setup: given the neighborhoods of two samples of a generative model, determine if the samples are identical, or are drawn independently. Part of the difficulty of the problem in this setup is that it may be required to determine if two neighborhoods are isomorphic or not. While we leave the question of graph isomorphism for randomly generated graphs for future work, we note that some of the techniques used for the classical graph isomorphism problem are related to our results. In particular our techniques for studying dense random graphs in Section 4.2 resemble some of the algorithms suggested for graph isomorphisms for some subclasses of graphs [Cai et al., 1992].

We also note that the question of whether an infinite graph is determined by some collection of its finite subgraphs has been studied in the context of unimodular and transitive infinite graphs [Aldous and Lyons, 2007] [Frisch and Tamuz, 2014].

1.1 General setup, models, main results

A (deterministic or random) graph $\mathcal{G} = \mathcal{G}_N$ with N vertices and labels (again possibly random) from a finite set on each vertex or edge is given. Each vertex v has a neighborhood $\mathcal{N}_r(v)$ of “radius” r which could be all of the vertices at distance r or some variation (see the examples below); we assume that location of vertex v is given in $\mathcal{N}_r(v)$.

- Q1. (Identifiability) Given each of the N neighborhoods $\mathcal{N}_r(v)$ for v a vertex in the network, can we correctly identify (up to isomorphism) the graph \mathcal{G} and its labels? We view this question as having two parts: (a) combinatorial criteria for identifiability (or non-identifiability), and (b) the probability of identifiability under particular random generative models.
- Q2. (Reconstruction) Assuming identifiability for a given \mathcal{G}_N and r , for $0 < \varepsilon < 1$, what is the minimum number, $M_{\text{rec}}(N, r, \varepsilon)$, of samples (with replacement) from the collection of neighborhoods that is necessary to ensure that the chance of correctly reconstructing the network \mathcal{G} with labels from the sample is at least $1 - \varepsilon$?

Questions Q1(a) and Q2 are discussed in Section 2, where we derive general results about combinatorial criteria for identifiability and upper and lower bounds on $M_{\text{rec}}(N, r, \varepsilon)$ based on coupon collecting. Notably our conditions for non-identifiability require that the graph is not isomorphic to small perturbations of the graph obtained by replacing a neighborhood with a non isomorphic neighborhood (thus avoiding the difficulty of the reconstruction conjecture). In Sections 3, 4, and 5,

Question Q1(b) is discussed and the general results of Section 2 are applied in the following three examples. Let $d(v, w)$ denote the distance between two vertices in a graph.

1. \mathcal{G} is the d dimensional n -lattice, here denoted \mathbf{Z}_n^d , with i.i.d. vertex labels from a probability distribution on $\{1, \dots, q\}$ and the neighborhoods $\mathcal{N}_r(v)$ are the $(n - r - 1)^d$ r -cubes. Here our neighborhoods differ slightly from the general setup and $N = N_{n,d,r} := (n - r - 1)^d$.
2. \mathcal{G} is an Erdős-Rényi random graph with vertex set V of size N and edge probability p_N where the vertices have no labels (or you can think of each having the same label) and the r -neighborhoods, $\mathcal{N}_r(v), v \in V$, are the subgraphs induced by the vertices at distance no greater than r from each vertex. We also consider labeled generalizations of the model.
3. *The random jigsaw puzzle problem.* \mathcal{G} is the $n \times n$ lattice and we view each vertex as being the center of a puzzle piece with each of the four edges receiving one of q jigs. Thus each vertex is labeled with an ordered 4-tuple of the q possible labels (jigs), corresponding to the label of each edge. Note that adjacent vertices have dependent labels. The neighborhoods $\mathcal{N}_0(v)$ are simply the vertices with labels and correspond to the puzzle pieces.

The main question we address in these examples is what are conditions on r or q as $N \rightarrow \infty$ to ensure identifiability (or non-identifiability)? We now summarize a subset of our findings and open problems.

Example 1: Lattices In Section 3, we find that if the vertices of the lattice are labeled uniformly and independently then, up to constants, the asymptotic threshold of r for identifiability is $\log(n)^{1/d}$.

Theorem 1.1. *For \mathbf{Z}_n^d with vertex labels i.i.d. uniform from fixed q labels and taking limits as $n \rightarrow \infty$, if for some $\varepsilon > 0$,*

$$r^d \leq (1 - \varepsilon) \frac{d}{2^{d-1}} \frac{\log n}{\log q},$$

then the probability of identifiability from r -neighborhoods tends to zero, and if for some $\varepsilon > 0$,

$$r^d \geq (1 + \varepsilon) 2d \frac{\log n}{\log q},$$

then the probability of identifiability from r -neighborhoods tends to one.

We conjecture that:

Conjecture 1.2. *There exists a constant $c_d(q)$ such that for every $\varepsilon > 0$, when $r^d \geq (1 + \varepsilon)c_d(q) \log n$, the probability of identifiability goes to 1 as $n \rightarrow \infty$, while when $r^d \leq (1 - \varepsilon)c_d(q) \log n$, the probability of identifiability goes to 0.*

More ambitiously we can ask:

Question 1.3. *Does there exist a constant c_d such that for every $\varepsilon > 0$, when $r^d \geq (1 + \varepsilon)c_d \frac{\log n}{\log q}$, the probability of identifiability goes to 1 as $n \rightarrow \infty$, while when $r^d \leq (1 - \varepsilon)c_d \frac{\log n}{\log q}$, the probability of identifiability goes to 0?*

In both cases finding the value of the constant, $c_d(q)$ or c_d , is a challenging open problem. The case of non-uniform labels is also discussed in Section 3.

Example 2: Erdős-Rényi graphs. The results of Section 4 show that for $\lambda \neq 1$, the asymptotic threshold for identifiability in the sparse Erdős-Rényi random graph is $\log(N)$ (up to constants).

Theorem 1.4. *For the Erdős-Rényi graph on N vertices with $p_N = \lambda/N$ and taking limits as $N \rightarrow \infty$, if for some $\varepsilon > 0$*

$$\frac{r}{\log(N)} < \frac{1}{2(\lambda - \log(\lambda))} - \varepsilon,$$

then the probability of identifiability from r -neighborhoods tends to zero.

- *If $\lambda < 1$ and for some $\varepsilon > 0$,*

$$\frac{r}{\log(N)} > \frac{1}{\log(1/\lambda)} + \varepsilon,$$

then the probability of identifiability from r -neighborhoods tends to one.

- *If $\lambda > 1$ and $\lambda_* < 1$ is the unique solution to $\lambda e^{-\lambda} = \lambda_* e^{-\lambda_*}$, and for some $\varepsilon > 0$,*

$$\frac{r}{\log(N)} > \frac{1}{\log(\lambda)} + \frac{2}{\log(1/\lambda_*)} + \varepsilon,$$

then the probability of identifiability from r -neighborhoods tends to one.

For $\lambda = 1$, the second statement of Theorem 4.2 below implies that if $rN^{-1/3} \rightarrow \infty$, then the probability of identifiability tends to one, but this is far from the lower bound $\log(N)$ provided by the previous result. We make the following conjecture:

Conjecture 1.5. *For positive $\lambda \neq 1$, there exists a constant c_λ such that for every $\varepsilon > 0$, when $r \geq (1+\varepsilon)c_\lambda \log N$, the probability of identifiability tends to 1 as $N \rightarrow \infty$, while when $r \leq (1-\varepsilon)c_\lambda \log N$, the probability of identifiability goes to 0.*

Natural open problems are to prove the conjecture, find the value of c_λ , and also to better understand the critical case where $\lambda = 1$. The cases of sparse Erdős-Rényi with labels and Erdős-Rényi with unbounded degrees are also studied in Section 4. In particular, in the most technical result in the paper we show that if $p_N = \omega(\log(N)^2/N)$ then neighborhoods of size 3 are enough to ensure identifiability:

Theorem 1.6. *If \mathcal{G} is the Erdős-Rényi random graph with N vertices and edge probability p_N satisfying $Np_N/\log(N)^2 \rightarrow \infty$ as $N \rightarrow \infty$ and we are given $\mathcal{N}_3(v)$ for each vertex v in \mathcal{G} , then the probability of identifiability tends to one.*

Example 3: Jigsaw puzzle. In Propositions 5.1 and 5.2, we show that if $q = o(n^{2/3})$, then the probability of identifiability tends to zero and if $q = \omega(n^2)$, then the probability of identifiability tends to one. We do not believe that either the constant $2/3$ or the constant 2 is sharp but conjecture there is a critical exponent:

Conjecture 1.7. *For the jigsaw puzzle problem, there exists a constant c such that for all $\varepsilon > 0$ if*

- *$q \leq n^{c-\varepsilon}$ then the probability of identification goes to 0 as $n \rightarrow \infty$ and if*
- *$q \geq n^{c+\varepsilon}$ then the probability of identification goes to 1 as $n \rightarrow \infty$.*

A number of additional open problems and conjectures are given in each section and we conclude the paper with a summary of these and other outstanding questions in Section 6. In this paper we only consider either unlabeled graphs or graphs that have i.i.d. labels. However, we emphasize that the questions considered here can be naturally extended to labelings of the graph outside of the i.i.d. case. For example the graph may be colored by an Ising model or by a uniform proper coloring. Thus the study of graph shotgun assembly raises new problems in random graphs, percolation, Ising/Potts models, as well as algorithmic problems regarding random constraint satisfaction problems and the theory of spin glasses.

Except for the case of dense ER random graphs and the DNA shotgun assembly problem, none of the graph shotgun results are tight. We conclude the introduction with another family of examples for which it is easy to derive tight bounds.

The labelled full binary tree. Let \mathcal{T}_n be the full binary tree with 2^n leaves and label each vertex uniformly from the letters $\{1, \dots, q\}$. We are given the 1-neighborhoods $\mathcal{N}_1(v)$ of the $2^n - 2$ vertices that are not leaves or the root (so we see the labels of the vertex, its two children, and its parent).

Proposition 1.8. *Let $\varepsilon > 0$. If*

$$\frac{\log(q)}{n} < \log(2) - \varepsilon,$$

then the probability of identifiability of the labeled binary tree \mathcal{T}_n from 1-neighborhoods tends to zero. If

$$\frac{\log(q)}{n} > \log(2) + \varepsilon,$$

then the probability of identifiability of the labeled binary tree \mathcal{T}_n from 1-neighborhoods tends to one.

Proof. To prove the first assertion, note that if there are two vertex disjoint edges between levels $n - 2$ and $n - 1$ of the tree having endpoints with identical labels, then with good probability reconstruction is impossible since we can switch the cherries below these edges (which have different labels with good probability) and obtain a non-isomorphic (again with good probability) labeling of the tree with the same neighborhoods. Thus we lower bound the probability of this event using the second moment method. Actually it's enough to consider neighborhoods of vertices at level $n - 1$ which are odd-numbered when labeled sequentially $1, 2, \dots, 2^{n-1}$ starting from the left. Let $B = B_{n,q}$ be the number of pairs of such neighborhoods where the central vertices have the same label and the parent vertices have the same label (possibly different from the central vertices) and the two pairs of leaves have different labels (as sets). Writing $B = \sum_{\alpha \neq \beta} X_{\alpha,\beta}$, where the sum is over all such pairs of neighborhoods (α, β) and $X_{\alpha,\beta}$ is the indicator of the event just described, we compute

$$\mathbb{E}B \leq 2^{2(n-2)}(1/q)^2(1 - 1/q^3).$$

After noting that the labels being chosen uniformly implies the $X_{\alpha,\beta}$ are independent, we find

$$\text{Var}B = \sum_{\alpha \neq \beta} \text{Var}(X_{\alpha,\beta}) \leq \mathbb{E}B,$$

and the first claim of the proposition follows by the second moment method.

For the second part of the claim, it's clear that if no two edges have the same labels, then we can piece together the tree from the neighborhoods by overlapping distinct edges. The mean of the number of pairs of edges with the same labels is bounded above by

$$2^{2n+2}(1/q)^2,$$

which tends to zero under the hypothesis of the second statement of the proposition and so the result follows. \square

2 Combinatorial and sampling results

We introduce two concepts that can be used to determine identifiability: blocking configurations and uniqueness of overlaps. For concreteness, specialize to the case where for each vertex v , $\mathcal{N}_r(v)$ is the labeled subgraph induced by the vertices at distance no greater than r from each vertex.

2.1 Blocking configurations

A blocking configuration is a neighborhood structure or pattern such that if it appears then identifiability is impossible. For a given example, there can be a number of different blocking configurations, though that described in Lemma 2.1 below is most likely in our examples. In random models, we use blocking configurations to get upper bounds on the asymptotic neighborhood size to ensure non-identifiability: if the neighborhoods grow too slowly, then the chance that a blocking configuration appears tends to one and identifiability is impossible (or the probability is bounded away from zero and so identifiability isn't assured).

For $t > s > 0$ and vertex v , define the sphere (or shell) $\mathcal{S}(v; s, t)$ to be subgraph induced by *edges* connecting vertices having distance to v between s and t (inclusive). Note that $\mathcal{S}(v; s, t)$ has no isolated vertices.

Lemma 2.1. *If \mathcal{G} is such that there is an $r > 0$ and vertices v, w such that*

- (i) $\mathcal{S}(v; 1, 2r) = \mathcal{S}(w; 1, 2r)$,
 - (ii) $d(v, w) > 2r$, and
 - (iii) *the graph obtained by switching $\mathcal{N}_1(v)$ and $\mathcal{N}_1(w)$ in \mathcal{G} is not isomorphic to \mathcal{G} ,*
- then identifiability from r -neighborhoods is impossible.*

Proof. We claim that there are at least two non-isomorphic labeled graphs having the same r -neighborhoods as \mathcal{G} : the true one, \mathcal{G} , and one where $\mathcal{N}_1(v)$ and $\mathcal{N}_1(w)$ are switched, denoted by \mathcal{G}' . Condition (i) ensures that such a switch is possible since the number of vertices at distance one connecting to vertices at distance two and their labels agree for v and w . Condition (iii) ensures that \mathcal{G} and \mathcal{G}' are not isomorphic (and note in particular that this implies $\mathcal{N}_1(v) \neq \mathcal{N}_1(w)$). Denote by \mathcal{N}'_r the r -neighborhoods generated by \mathcal{G}' .

We only need to show that \mathcal{G} and \mathcal{G}' generate the same r -neighborhoods (including multiplicities). From (ii), there is no vertex having both v and w in its \mathcal{G} r -neighborhood. Thus we can split vertices into two groups: those being within distance r of exactly one of v or w in \mathcal{G} , and those having distance greater than r from both of v and w . For any vertex x in the latter group, the differences in switching $\mathcal{N}_1(v)$ and $\mathcal{N}_1(w)$ are not reflected by (potential) neighbors of v and w that are at distance r from x (since the labels and positions of such vertices have to match), and so $\mathcal{N}_r(x) = \mathcal{N}'_r(x)$.

For the group of vertices within distance r (in \mathcal{G}) of one of v or w , Condition (i) implies there is an obvious matching of each vertex x that satisfies either

- $2 \leq d(x, v) \leq r$ (distance in \mathcal{G}) or,
- $d(x, v) = 1$ and x has a neighbor at distance two from v ,

to one having the same distance from w and identical label. Moreover, under this matching, $\mathcal{N}_r(x) = \mathcal{N}'_r(y)$ and $\mathcal{N}'_r(x) = \mathcal{N}_r(y)$. Finally, by (i), for $x = v, w$ or a neighbor of v or w with no neighbors at distance 2 from v or w , $\mathcal{N}_r(x) = \mathcal{N}'_r(x)$. Thus \mathcal{G} and \mathcal{G}' generate the same r -neighborhoods. \square

Remark 2.2. Condition (iii) seems a bit unnatural and possibly hard to verify. Indeed, it is difficult to check in situations where the graph G has many symmetries since the graph isomorphism problem is computationally difficult. However, such symmetry is rare in random graphs and so in our applications of the lemma, Condition (iii) is easy to verify. We also note that the condition is reasonable to impose given the difficulty of the “reconstruction conjecture” that has been open for more than 50 years.

2.2 Uniqueness of overlaps

The next result formalizes the intuition that if all of the neighborhoods of a certain size are unique, then slightly larger neighborhoods are enough to ensure identifiability. In random models, we use uniqueness of overlaps to get lower bounds on the asymptotic neighborhood size to ensure identifiability. If the neighborhoods grow quickly enough, then the chance that all neighborhoods of a slightly smaller size are unique tends to one and identifiability is ensured.

Lemma 2.3. *If $\mathcal{N}_{r-1}(v) \neq \mathcal{N}_{r-1}(w)$ for all vertices $v \neq w$, then there is an efficient algorithm for recovering the graph from r -neighborhoods.*

Proof. We can sequentially build the network by overlapping neighborhoods of radius $r - 1$. Start with some r -neighborhood $\mathcal{N}_r(v)$ and note that the $(r - 1)$ -neighborhood of each neighbor of v is contained in $\mathcal{N}_r(v)$ and these are all unique by assumption. Thus for each vertex $w \neq v$, we examine the $(r - 1)$ -neighborhoods of neighbors of w and overlap any of these matching the $(r - 1)$ -neighborhoods of neighbors of v . Repeating this process for each neighbor of v and then continuing for the vertices at distance 2, 3, \dots from v , it’s clear that the process terminates when a connected component is recovered. \square

Remark 2.4. The proof of the lemma is simple because we assume we see not only $\mathcal{N}_r(v)$, but also which vertex in the neighborhood is the “center” (namely, v). We do not investigate here how to relax this condition to the situation where the center v is not given.

2.3 Sampling

In the regime where we have uniqueness of $(r - 1)$ -neighborhoods, then the coupon collector problem yields bounds on the probability of reconstruction. Let $M_{\text{rec}}(N, r, \varepsilon)$ be the minimum number of samples so that the chance the graph can be reconstructed from the samples is least $1 - \varepsilon$.

Lemma 2.5. *If for some r , $\mathcal{N}_{r-1}(v) \neq \mathcal{N}_{r-1}(w)$ for all vertices $v \neq w$, then*

$$M_{\text{rec}}(N, r, \varepsilon) \leq \lceil N \log(N) - N \log \varepsilon \rceil.$$

Proof. The proof of Lemma 2.3 implies that it’s enough to see all of the neighborhoods, possibly in multiplicities, since then we can build the network by overlapping the $(r - 1)$ -neighborhoods of neighbors of the sampled vertex. The bound in the lemma now easily follows from coupon collecting: if T is the number of samples with replacement required to collect N distinct coupons, then a union bound implies that for integer $M > 0$,

$$\mathbb{P}(T > M) \leq N(1 - 1/N)^M \leq Ne^{-M/N}.$$

Now setting $M = \lceil N \log(N) - N \log \varepsilon \rceil$, we find

$$\mathbb{P}(\text{Can't reconstruct with } M \text{ samples}) \leq \mathbb{P}(T > M) \leq \varepsilon,$$

and so $M_{\text{rec}}(N, r, \varepsilon) \leq M$. \square

Since there is no hope of reconstruction if there is some vertex that doesn't appear in any of the sampled neighborhoods, we can also use coupon collecting to get a lower bound on $M_{\text{rec}}(N, r, \varepsilon)$ in the general case. Let $|\mathcal{N}_r(v)|$ denote the number of vertices in $\mathcal{N}_r(v)$.

Lemma 2.6. *If M is such that*

$$\frac{\left(\sum_{i=1}^N \left(1 - \frac{|\mathcal{N}_r(v_i)|}{N} \right)^M \right)^2}{\sum_{i,j=1}^N \left(1 - \frac{|\mathcal{N}_r(v_i) \cup \mathcal{N}_r(v_j)|}{N} \right)^M} \geq \varepsilon,$$

then $M_{\text{rec}}(N, r, \varepsilon) \geq \lfloor x \rfloor$.

Proof. Let W_M be the number of vertices that have not appeared in some neighborhood in a sample of size M . If $W_M > 0$, then we can't reconstruct with M samples and so by the second moment method,

$$\mathbb{P}(\text{Can't reconstruct with } M \text{ samples}) \geq \mathbb{P}(W_M > 0) \geq \frac{(\mathbb{E}W_M)^2}{\mathbb{E}W_M^2}, \quad (2.1)$$

and for any M such that the right-most side of (2.1) is greater than ε , the chance of reconstruction is at most $1 - \varepsilon$ which implies $M \leq M_{\text{rec}}(N, r, \varepsilon)$. The result now follows by computing

$$\begin{aligned} \mathbb{E}W_M &= \sum_{i=1}^N \left(1 - \frac{|\mathcal{N}_r(v_i)|}{N} \right)^M, \\ \mathbb{E}W_M^2 &= \sum_{i,j=1}^N \left(1 - \frac{|\mathcal{N}_r(v_i) \cup \mathcal{N}_r(v_j)|}{N} \right)^M. \end{aligned} \quad \square$$

3 Labeled lattice models

Recall the setting of Example 1: \mathcal{G} is the $d \geq 2$ dimensional n -box \mathbf{Z}_n^d with i.i.d. vertex labels and neighborhoods the r -boxes contained in \mathbf{Z}_n^d ; note that for these neighborhoods the position of v is irrelevant. Our results for i.i.d. uniform labeling are different than the general i.i.d. case.

3.1 Uniform labels

Assume the vertices of \mathbf{Z}_n^d are labeled uniformly from $q \geq 2$ labels. Our first result uses blocking configurations to obtain an upper bound on the growth of r to ensure a positive chance of non-identifiability.

Proposition 3.1. *Given the r -neighborhoods of \mathbf{Z}_n^d with vertex labels i.i.d. uniform from q labels, the following holds as $n \rightarrow \infty$.*

- *if $(n/r)^{2d} q^{-(2r)^d} \rightarrow \infty$, then the probability of identifiability tends to zero, and*
- *if $\liminf_{n \rightarrow \infty} \left[(n/r)^{2d} q^{-(2r)^d} \right] > 0$, then the probability of identifiability is strictly less than one.*

Proof. We lower bound the probability of the following blocking configurations given by a pair of *non-overlapping* $(2r-1)$ -neighborhoods that have identical labels except for the two center vertices which are different. We only consider neighborhoods of the form $x + [0, 2r-1]^d$ where all of the coordinates of x are 0 modulo $2r-1$. Similar to Lemma 2.1, if two such neighborhoods exist, then identifiability is impossible since there are at least two ways to construct a consistent layout of neighborhoods, by switching the labels of the center vertices. Note further that the probability that there is an isomorphism of the graph excluding these two neighborhoods is at most $2^d \times (1/q)^{n^d/2-1}$ (since there are 2^d possible rotations and each site has to match the label of one other site).

To establish the existence of the neighborhood pair, we use the second moment method. Let $B = B_{n,d,r,q}$ denote the number of such blocking configurations described above and we compute $\mathbb{E}B$ and $\mathbb{E}B^2$. Assume that $n \gg r$ (without loss under the hypotheses of the proposition) and denote the set of such $(2r-1)$ -neighborhoods of \mathbf{Z}_n^d by $\Gamma = \Gamma_{n,d,2r-1}$ (note that $|\Gamma| = \Theta((n/2r)^d)$) and write

$$B = \sum_{\alpha, \beta \in \Gamma, \alpha \cap \beta = \emptyset} X_{\alpha, \beta},$$

where $X_{\alpha, \beta}$ is the indicator of the event that the labels of α and β are equal except for the center labels which must be different and $\alpha \cap \beta = \emptyset$ means α and β are non-overlapping. It's easy to see that $\mathbb{E}X_{\alpha, \beta} = (1/q)^{(2r-1)^d-1} (1 - 1/q)$ which implies that

$$\mathbb{E}B \geq \Theta((n/2r)^{2d})(1/q)^{(2r-1)^d-1} (1 - 1/q), \quad (3.1)$$

The fact that B is concentrated follows from the fact that the $X_{(\alpha, \beta)}$ are pairwise independent: *if the labels are chosen uniformly*, then for two pairs of neighborhoods $(\alpha, \beta) \neq (\gamma, \delta)$, $X_{\alpha, \beta}$ and $X_{\gamma, \delta}$ are independent. Thus

$$\text{Var}(B) = \sum_{\alpha, \beta \in \Gamma, \alpha \cap \beta = \emptyset} \text{Var}(X_{\alpha, \beta}) \leq \mathbb{E}B.$$

Now the proof follows by the second moment method. □

We can use uniqueness of overlaps as in Lemma 2.3 to find a regime where asymptotic reconstruction is assured.

Proposition 3.2. *If $n^{2d}q^{-(r-1)^d} \rightarrow 0$ as $n \rightarrow \infty$, then the probability of identifiability (of \mathbf{Z}_n^d with i.i.d. uniform on q vertex labels) from r -neighborhoods tends to one.*

Proof. Let $Y := Y_{n,d,r,q}$ be the number of pairs of different $(r-1)$ -neighborhoods that have the same labels and we show that $\mathbb{E}Y \rightarrow 0$ as $n \rightarrow \infty$, from which the result follows from a minor variation of Lemma 2.3.

Similar to the proof of Proposition 3.1, denote the set of $(r-1)$ -neighborhoods of \mathbf{Z}_n^d by $\Gamma = \Gamma_{n,d,r-1}$ and for $\alpha, \beta \in \Gamma$, let $Y_{(\alpha, \beta)}$ be the indicator that α and β have the same labels. It's obvious that if $\alpha \cap \beta = \emptyset$ (meaning the two neighborhoods share no vertices) then $\mathbb{E}Y_{(\alpha, \beta)} = q^{-(r-1)^d}$, but *since the labels are uniform*, straightforward considerations (see below) show that in fact

$$\mathbb{E}Y_{(\alpha, \beta)} = q^{-(r-1)^d} \quad (3.2)$$

for all $\alpha \neq \beta$. Thus we find

$$\mathbb{E}Y = \sum_{\alpha, \beta \in \Gamma, \alpha \neq \beta} \mathbb{E}Y_{(\alpha, \beta)} = [(n-r)^{2d} - 1]q^{-(r-1)^d}.$$

To prove (3.2) formally assume WLOG that $(x, y) \rightarrow (x, y) - (i, j)$ is an injective map from β to α where $i, j \geq 0$ and at least one of i and j is non-zero. Then we can label $\alpha \cup \beta$ according to lexicographic order where

- If a site is in $\alpha \setminus \beta$ then we label it arbitrarily.
- If it is in β then we label it by looking at the site $(x, y) - (i, j)$ which was already labeled.

This defines all labelings of $\alpha \cup \beta$ where α and β have the same label so the number of such labelings is $q^{|\alpha \setminus \beta|}$ while the total number of labelings of $\alpha \cup \beta$ is $q^{|\alpha \cup \beta|}$. The proof follows. \square

Theorem 1.1 in the introduction is easily established by combining Propositions 3.1 and 3.2.

3.2 Non-uniform labels

If the labels are i.i.d. but not uniform, we can prove a (weaker) analog of Proposition 3.2. Let p_i denote the chance of label i appearing at a site and $\mathcal{P}_j = \sum_i p_i^j$ denote the probability that j particular sites have the same label.

Proposition 3.3. *If $(nr)^{2d} \mathcal{P}_2^{(r-1)^d} \rightarrow 0$ as $n \rightarrow \infty$, then the probability of identifiability (of \mathbf{Z}_n^d with i.i.d. vertex labels) from r -neighborhoods tends to one.*

Proof. As in the proof of Proposition 3.2, let Y be the number of $(r-1)$ -neighborhoods that have the same labels and we show that $\mathbb{E}Y \rightarrow 0$ as $n \rightarrow \infty$. Similar to the proof of Proposition 3.1, denote the set of $(r-1)$ -neighborhoods of \mathbf{Z}_n^d by $\Gamma = \Gamma_{n,d,r-1}$ and for $\alpha, \beta \in \Gamma$, let $Y_{(\alpha,\beta)}$ be the indicator that α and β have the same labels. It's obvious that if $\alpha \cap \beta = \emptyset$ (meaning the two neighborhoods share no vertices) then $\mathbb{E}Y_{(\alpha,\beta)} = \mathcal{P}_2^{(r-1)^d}$. If $\alpha \cap \beta \neq \emptyset$, then

$$\mathbb{E}Y_{(\alpha,\beta)} = \prod_{j \geq 2} \mathcal{P}_j^{k_j}, \quad (3.3)$$

where $j \times k_j$ are the number of sites in the union of α and β that need to be matched to $j-1$ other sites to ensure $Y_{(\alpha,\beta)} = 1$ (c.f., the justification of (3.2) at the end of the proof of Proposition 3.2). Note that $\sum_{j \geq 2} (j-1)k_j = (r-1)^d$ and that $\sum_{j \geq 2} k_j = |\alpha \cup \beta| - (r-1)^d$, since this sum is equal to $|\alpha/\beta|$. Using the basic inequality $\mathcal{P}_j \leq \mathcal{P}_2^{j/2}$ for $j \geq 2$ in (3.3), we find

$$\mathbb{E}Y_{(\alpha,\beta)} \leq \prod_{j \geq 2} \mathcal{P}_2^{jk_j/2} = \mathcal{P}_2^{|\alpha \cup \beta|/2} \leq \mathcal{P}_2^{(r-1)^d/2};$$

the last inequality is since $|\alpha \cup \beta| \geq (r-1)^d$. Counting the number of overlapping and non-overlapping neighborhoods, we find

$$\mathbb{E}Y \leq n^{2d} \mathcal{P}_2^{(r-1)^d} + 4r^d n^d \mathcal{P}_2^{(r-1)^d/2},$$

from which the result easily follows. \square

Remark 3.4. If the labels are uniform, then $\mathcal{P}_j = q^{-(j-1)}$ and so we can use this exact quantity (rather than the inequality $\mathcal{P}_j \leq \mathcal{P}_2^{j/2}$) in (3.3) in the proof of Proposition 3.3 to recover the sharper Proposition 3.2.

For non-uniform vertex labels, the correlations between the appearance of overlapping blocking sets can become significant and so the second moment method of Proposition 3.1 breaks down. Still we believe that similar results should hold:

Conjecture 3.5. *Consider a distribution π that is fully supported on $\{1, \dots, q\}$ and the labeling of \mathbf{Z}_n^d by i.i.d. labels from π . For every dimension d , there exists a constant $c_d(\pi)$ such that for every $\varepsilon > 0$, when $r^d \geq (1 + \varepsilon)c_d(\pi) \log n$, the probability of identifiability tends to one as $n \rightarrow \infty$, while when $r^d \leq (1 - \varepsilon)c_d(\pi) \log n$, the probability of identifiability goes to 0.*

We believe that conjecture 3.5 should also extend to some dependent setups including:

- The uniform distribution of legal vertex colorings of a box with $q \geq 3d$ colors. We require that q is large to ensure correlation decay of the distribution. Note for example that if $q = 2$ and $d \geq 2$, then the problem is degenerate as there are only two possible colorings of the graph.
- The Ising and Potts models with finite temperature $0 < \beta < \infty$ in the box.

Proving the conjectures and establishing the value of the threshold in these examples are fascinating open problems.

3.3 Sampling

If \mathbf{Z}_n^d has uniqueness of $(r - 1)$ -overlaps (asymptotically assured in the regimes of Propositions 3.2 and 3.3), then the argument of Lemma 2.5 automatically implies an upper bound of $N(\log(N) - \log(N))$ (recall $N = N_{n,d,r} := (n - r - 1)^d$ is the number of neighborhoods) on $M_{\text{rec}}(N, \varepsilon, r)$, the minimum number of samples needed to reconstruct the abels of the lattice with probability at least $1 - \varepsilon$. We can also use Lemma 2.6 to show that we need at least of order (large N , small ε) $\frac{N}{r^d} (\log(N/r^d) - (\log(\varepsilon)))$ samples to reconstruct in any regime.

Proposition 3.6. *For \mathbf{Z}_n^d with vertex labels,*

$$M_{\text{rec}}(N, \varepsilon, r) \geq \frac{\log\left(\frac{1}{\varepsilon} - 1\right) - \log\left(\frac{(2r-1)^d}{N}\right)}{-\log\left(1 - \frac{r^d}{N}\right)}. \quad (3.4)$$

Proof. We may use Lemma 2.6 with this neighborhood structure since its argument only relies on the size (and not the structure) of the neighborhoods. First note $|\mathcal{N}_r(v)| = r^d$ for all v and $|\mathcal{N}_r(v) \cup \mathcal{N}_r(w)| = 2r^d$ if $\mathcal{N}_r(v) \cap \mathcal{N}_r(w) = \emptyset$ and $|\mathcal{N}_r(v) \cup \mathcal{N}_r(w)| \geq r^d$ otherwise. Using these bounds, if M is no greater than the right hand side of (3.4), then

$$\begin{aligned} \frac{\left(\sum_{i=1}^N \left(1 - \frac{|\mathcal{N}_r(v_i)|}{N}\right)^M\right)^2}{\sum_{i,j=1}^N \left(1 - \frac{|\mathcal{N}_r(v_i) \cup \mathcal{N}_r(v_j)|}{N}\right)^M} &\geq \frac{N^2 \left(1 - \frac{r^d}{N}\right)^{2M}}{N^2 \left(1 - \frac{2r^d}{N}\right)^M + N(2r-1)^d \left(1 - \frac{r^d}{N}\right)^M} \\ &\geq \left[1 + \frac{(2r-1)^d}{N} \left(1 - \frac{r^d}{N}\right)^{-M}\right]^{-1} \geq \varepsilon, \end{aligned}$$

and the result follows. \square

4 Erdős-Rényi graph

Assume the setup of Example 2: \mathcal{G} is the Erdős-Rényi random graph with N vertices, the vertices have no labels (or to fit our setup, all labels are the same) and for each vertex v , we have the r -neighborhoods $\mathcal{N}_r(v)$ which are the subgraphs induced by vertices at distance $\leq r$ from v . This example fits exactly into our general setup and so Lemmas 2.1 and 2.3 can be applied “out of the box”. As is typical for Erdős-Rényi random graphs, the results differ if the graph has bounded average degree or not and so we separate our results accordingly to Sections 4.1 and 4.2.

4.1 Bounded average degree Erdős-Rényi

Let \mathcal{G} be the Erdős-Rényi random graph with N vertices and edge probability $p_N = \lambda/N$ for some $\lambda > 0$. We use the blocking configuration of Lemma 2.1 to show the following result.

Proposition 4.1. *For the Erdős-Rényi graph on N vertices with $p_N = \lambda/N$, using the notation of the previous paragraph, and taking limits as $N \rightarrow \infty$,*

- *if $\sqrt{N}\lambda^r(1 - \lambda/N)^{Nr} \rightarrow \infty$, then the probability of identifiability tends to zero, and*
- *if $\liminf_{N \rightarrow \infty} \sqrt{N}\lambda^r(1 - \lambda/N)^{Nr} > 0$, then the probability of identifiability is strictly less than one.*

Proof. Note that $\lambda(1 - \lambda/N)^N < 1$ and so if r grows faster than $\log(N)$, then neither of the hypotheses of the proposition are satisfied, and so we can assume without loss that $r/N^a \rightarrow 0$ for all $a > 0$. We lower bound the probability of the appearance of the following blocking (induced) subgraph on $4r + 6$ vertices: the subgraph has two components, one a line graph on $2r + 1$ vertices and the other a line graph on $2r + 1$ vertices with the addition of both end vertices being connected to two other vertices with no other edges to form “prongs”; see Figure 1.

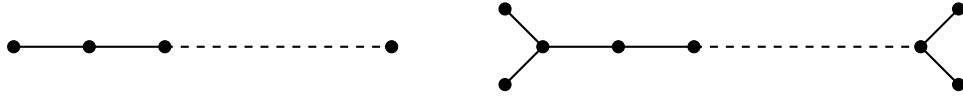


Figure 1: Example of blocking subgraph for neighborhoods of radius r . The line graph has $2r + 1$ vertices.

Note that this blocking set satisfies the hypotheses of Lemma 2.1 by taking v to be an endpoint of the line graph and w to be one of the degree three vertices. Alternatively, it's easy to see that if such a subgraph is present, then identifiability is impossible because there are at least two ways to construct the graph consistent with the neighborhoods, by switching one of the prongs to the line graph; see Figure 2 for illustration.

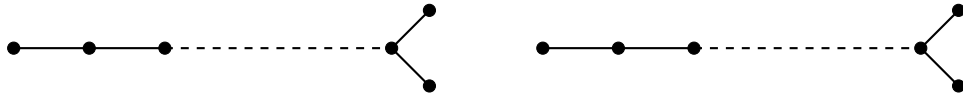


Figure 2: A subgraph that has the same r -neighborhoods as that of Figure 1

Let $B = B_{N,r,\lambda}$ be the number of such (induced) subgraphs of \mathcal{G} and write $B = \sum_{\alpha \in \Gamma} X_\alpha$, where $\Gamma = \Gamma_{N,4r+6}$ is the collection of subsets of vertices of size $4r + 6$ and for $\alpha \in \Gamma$, X_α is the indicator that the blocking subgraph of Figure 1 is the induced subgraph of \mathcal{G} on α . The X_α are equally

distributed and for $\alpha \neq \beta$, if $\alpha \cap \beta \neq \emptyset$, then $X_\alpha X_\beta = 0$. Thus we find for (say) $\alpha = \{1, \dots, 4r+6\}$ and $\beta = \{4r+7, \dots, 4r+12\}$,

$$\mathbb{E}B = \binom{N}{4r+6} \mathbb{E}X_\alpha, \quad \mathbb{E}B^2 = \mathbb{E}B \left(1 + \binom{N-4r-6}{4r+6} \mathbb{E}[X_\beta | X_\alpha = 1] \right).$$

From this point we need to compute $\mathbb{E}X_\alpha$ and $\mathbb{E}[X_\beta | X_\alpha = 1]$. There is at most one copy of the blocking (induced) subgraph on α , but there are a number of ways the subgraph can appear. By enumeration and noting the chance that any potential way the subgraph can appear, we find

$$\mathbb{E}X_\alpha = \binom{4r+6}{2r+1} \binom{2r+5}{4} \binom{4}{2} \frac{2(2r+1)!^2}{2^2} p_N^{2(2r+2)} (1-p_N)^{(4r+6)(N-3)+4}; \quad (4.1)$$

the first binomial coefficient counts the number of ways of assigning $2r+1$ vertices of α to the line graph, the second assigns four of the remaining vertices to the prongs and for each of the $(2r+1)$ -lines, there are $(2r+1)!/2$ ways to put them in order; the final factor of 2 comes from assigning the pairs of prong vertices to an end. Once the vertices are assigned, there are $4r+4$ edges that must appear, each with probability p_N , and $2(2r-1)(N-3) + 6(N-2) + 2(N-4)$ edges that must not appear.

Similarly, given $X_\alpha = 1$, none of the vertices of α have edges connecting to vertices outside of α and so $X_\beta | X_\alpha = 1$ is distributed as X_β , but on an Erdős-Rényi graph on $N-4r-6$ vertices and chance of edge p_N . Thus we use (4.1) but with $N-4r-6$ replacing N (except in p_N) to find

$$\mathbb{E}[X_\beta | X_\alpha = 1] = \binom{4r+6}{2r+1} \binom{2r+5}{4} \binom{4}{2} \frac{2(2r+1)!^2}{2^2} p_N^{2(2r+2)} (1-p_N)^{(4r+6)(N-2r+3)+4}. \quad (4.2)$$

Putting together (4.1) and (4.2) and using that under either of the hypotheses of the proposition, $r/N^a \rightarrow 0$ for any $a > 0$, we find

$$\frac{(\mathbb{E}B)^2}{\mathbb{E}B^2} \geq \frac{(N-4r-6)^{4r+6} p_N^{4r+4} (1-p_N)^{N(4r+6)}}{8 + (N-4r-6)^{4r+6} p_N^{4r+4} (1-p_N)^{(N-2r)(4r+6)}},$$

and under the first hypothesis of the proposition, the numerator and the denominator tend to infinity at the same rate, and under the second, the numerator on the right hand side stays bounded away from zero. \square

If r is larger than the diameter of the graph, then clearly we can identify from the neighborhoods. Thus we can use known results on the diameter of the Erdős-Rényi random graph (see [Riordan and Wormald, 2010], [Łuczak, 1998], [Nachmias and Peres, 2008], [Addario-Berry et al., 2012]) to get a lower bound on the growth of r to guarantee identifiability. Denote convergence in probability by \xrightarrow{p} .

Theorem 4.2. *Let \mathcal{G}_N be the Erdős-Rényi random graph on N vertices with edge probability $p_N = \lambda/N$ for a fixed $\lambda > 0$ and let $D = D_{N,\lambda}$ to be the maximum diameter of a component of \mathcal{G}_N .*

- [Łuczak, 1998, Theorem 11] *If $\lambda < 1$, then $D_{N,\lambda}/\log(N) \xrightarrow{p} 1/\log(1/\lambda)$.*
- [Nachmias and Peres, 2008, Theorem 1.1], [Addario-Berry et al., 2012, Theorem 5] *If $\lambda = 1$, then $N^{-1/3}D_{N,1}$ converges in distribution to a non-negative and non-degenerate distribution.*

- [Riordan and Wormald, 2010, Theorem 1.1] If $\lambda > 1$, and $\lambda_* < 1$ is the unique solution to $\lambda e^{-\lambda} = \lambda_* e^{-\lambda_*}$, then

$$\frac{D_{N,\lambda}}{\log(N)} \xrightarrow{p} \frac{1}{\log(\lambda)} + \frac{2}{\log(1/\lambda_*)}.$$

Theorem 1.4 in the introduction summarizes the lower bound on the neighborhood size for identifiability given by Proposition 4.1 and the upper bounds given by the properties of the diameter of Theorem 4.2.

Labeled Erdős-Rényi. Assuming vertices have i.i.d. labels from a finite set and we let \mathcal{P}_2 be the chance that two given vertices have the same label, we show the following result.

Proposition 4.3. *For the labeled Erdős-Rényi graph with $p_N = \lambda/N$, using the notation of the previous paragraph, and assuming $\mathcal{P}_2 \neq 1$, if for some $\varepsilon > 0$,*

$$\frac{r}{\log(N) - 2\log(1 - \mathcal{P}_2)} < \frac{1}{2\lambda - 2\log(\lambda) - \log(\mathcal{P}_2)} - \varepsilon,$$

then the chance of identifiability tends to zero as $N \rightarrow \infty$.

Proof. The argument is nearly identical to the proof of Proposition 4.1 but now the blocking configuration is two isolated line graphs with $2r + 1$ vertices, both having the same labels in the $2r - 1$ middle vertices, and each having two different labels at the endpoints; switching labels of an appropriately chosen endpoint (being careful of symmetries) from each line graph results in a non-isomorphic labeled graph with the same neighborhoods. If B is the number of such configurations, then the result follows from the second moment method after computing

$$\frac{(\mathbb{E}B)^2}{\mathbb{E}B^2} \geq \frac{(N - 4r - 2)^{4r+2} p_N^{4r} (1 - p_N)^{4Nr} \mathcal{P}_2^{2r-1} (1 - \mathcal{P}_2)^2}{8 + (N - 4r - 2)^{4r+2} p_N^{4r} (1 - p_N)^{(N-2r-4)(4r-2)} \mathcal{P}_2^{2r-1} (1 - \mathcal{P}_2)^2}. \quad \square$$

We make the following conjecture.

Conjecture 4.4. *Consider a distribution π that is fully supported on $\{1, \dots, q\}$ and the i.i.d. π -vertex labeling of the Erdős-Rényi random graph on N vertices with parameter λ/N . For positive $\lambda \neq 1$, there exists a constant $c_\lambda(\pi)$ such that for every $\varepsilon > 0$, when $r \geq (1 + \varepsilon)c_\lambda(\pi) \log N$, the probability of identifiability tends to one as $N \rightarrow \infty$, while when $r \leq (1 - \varepsilon)c_\lambda(\pi) \log N$, the probability of identifiability tends to 0.*

Open problems are to establish the conjecture, determine the value of $c_\lambda(\pi)$, and understand the critical case where $\lambda = 1$.

4.2 Dense Erdős-Rényi graph

Now we assume that \mathcal{G} is the Erdős-Rényi random graph with N vertices and edge probability p_N such that as $N \rightarrow \infty$, $Np_N/\log(N)^2 \rightarrow \infty$ and the neighborhoods are as before, described in Example 2. We restate and prove Theorem 1.6 from the introduction.

Theorem 4.5. *If \mathcal{G} is the Erdős-Rényi random graph with N vertices and edge probability p_N satisfying $Np_N/\log(N)^2 \rightarrow \infty$ as $N \rightarrow \infty$ and we are given $\mathcal{N}_3(v)$ for each vertex v in \mathcal{G} , then the probability of identifiability tends to one.*

Proof. If $p_N > N^{-3/5}$, then the diameter of \mathcal{G} is at most 3 [Bollobás, 1981] and so we can assume without loss that $p_N \leq N^{-3/5}$.

We show that the chance of the event “each vertex v has distinct 2-neighborhood $\mathcal{N}_2(v)$ ” tends to one and then the result follows by the uniqueness of overlaps Lemma 2.3. If v and w are distinct vertices of \mathcal{G} , then it’s enough to show as $N \rightarrow \infty$,

$$N^2 \mathbb{P}(\mathcal{N}_2(v) = \mathcal{N}_2(w)) \rightarrow 0. \quad (4.3)$$

In order for $\mathcal{N}_2(v) = \mathcal{N}_2(w)$, the degree of v ($\deg(v)$) must be equal to that of w and the degrees of the neighbors of v and w must be equal as multi-sets. Note that we can write $\deg(v) = B_v + I$ and $\deg(w) = B_w + I$ where B_v and B_w are independent with distribution $\text{Bi}(N-2, p_N)$ and I is the indicator that v and w have an edge between them. We bound the chance that v and w have the same degree and the chance of sharing too many neighbors as follows.

1. The Chernoff bound of Lemma 4.6 applied to the binomial distribution implies that for all $0 < \varepsilon_1 < 1/2$,

$$\mathbb{P}(\deg(v) \in Np_N(1 \pm \varepsilon_1)) \geq 1 - 2 \exp \left\{ -\frac{\varepsilon_1^2}{3} Np_N \right\}.$$

2. Noting that the event $\deg(v) = \deg(w)$ is independent of I , the indicator that v and w have an edge between them, we use the local limit theorem for the binomial distribution to find for C not depending on N ,

$$\mathbb{P}(\deg(w) = \deg(v) | \deg(v) \in Np_N(1 \pm \varepsilon_1)) \leq \frac{C}{\sqrt{Np_N}}.$$

3. Let $M := \deg(v) \mathbb{I}[\deg(w) = \deg(v) \in Np_N(1 \pm \varepsilon_1)]$ be the common degree of v and w assuming the conditioning of the items above hold (and zero otherwise), and let $K = M - |\mathcal{N}_1(v) \cap \mathcal{N}_1(w)| + I$ be the number of neighbors of v and w that are connected to exactly one of v or w . Given $M > 0$, the neighbors of v and w are each chosen uniformly from the $N-1$ possible neighbors. Thus if v and w are not neighbors, then $M-K$ is hypergeometric with M draws, M marked balls and $N-2$ total balls and if v and w are neighbors, then $M-K$ is hypergeometric with $M-1$ draws, $M-1$ marked balls and $N-2$ total balls. In either case, after noting that hypergeometric distributions can be represented as sums of *independent* indicators [Pitman, 1997] the Chernoff bound of Lemma 4.6 implies that for $0 < \varepsilon_2 < 1/2$,

$$\mathbb{P}(M-K \in M^2 N^{-1}(1 \pm \varepsilon_2) | M > 0) \geq 1 - 2 \exp \left\{ -\frac{\varepsilon_2^2}{3} Np_N \right\}.$$

Note that if $M > 0$, then $M \in Np_N(1 \pm \varepsilon_1)$ and so the event $M-K \in M^2 N^{-1}(1 \pm \varepsilon_2)$ implies $M-K \in Mp_N(1 \pm \varepsilon_1 \pm \varepsilon_2)$ and that

$$\begin{aligned} K &> M(1 - p_N(1 + \varepsilon_1 + \varepsilon_2)) > Np_N(1 - 2p_N)/2, \\ K &< M(1 - p_N(1 - \varepsilon_1 - \varepsilon_2)) < 2Np_N. \end{aligned}$$

Given $M > 0$ and K , let $\{D(v)\} := \{D_1(v), \dots, D_K(v)\}$ and $\{D(w)\} := \{D_1(w), \dots, D_K(w)\}$ denote the multi-set of degrees of the K non-intersecting neighbors of v and w , respectively. The three items above imply the following bound.

$$\mathbb{P}(\mathcal{N}_2(v) = \mathcal{N}_2(w)) \leq 2 \exp \left\{ -\frac{\varepsilon_1^2}{3} Np_N \right\} + \frac{2C \exp \left\{ -\frac{\varepsilon_2^2}{3} Np_N \right\}}{\sqrt{Np_N}} \quad (4.4)$$

$$+ \frac{C}{\sqrt{Np_N}} \mathbb{P}[\{D(v)\} = \{D(w)\} | M > 0, Np_N(1 - 2p_N)/2 < K < 2Np_N]. \quad (4.5)$$

Since $p_N \geq \log(N)^2/N$, the first two terms of (4.4) are easily seen to be $o(1/N^2)$ so we only need to bound (4.5).

Write $D_i(v) = V_i + A_i$, where V_i is the number of edges between v and the $N - 2K - 1$ vertices not in $(\mathcal{N}_1(v) \cup \mathcal{N}_1(w)) / (\mathcal{N}_1(v) \cap \mathcal{N}_1(w))$ and A_i the number edges between v and the remaining $2K - 1$ potential neighbors. Similarly, write $D_i(w) = W_i + B_i$. Note that given $M > 0$ and K , $\{V_1, \dots, V_K\}$ and $\{W_1, \dots, W_K\}$ are two independent (unordered) collections of i.i.d. $\text{Bi}(N - 2K - 1, p_N)$ random variables that are also independent of the A_i 's and B_i 's.

We show (i) that with good probability A_i and B_i are bounded by a constant and (ii) that the chance that independent binomial multi-sets are within constants is small.

For (i), the A_i 's and B_i 's form the degrees of an Erdős-Rényi graph on $2K$ vertices with edge parameter p_N and so are each marginally distributed $\text{Bi}(2K - 1, p_N)$. Thus,

$$\mathbb{P}(A_i < x, B_i < x, i = 1, \dots, K) \geq 1 - 2K\mathbb{P}(A_1 > x) \geq 1 - 2K \left(\frac{e}{x}\right)^x ((2K - 1)p_N)^x,$$

where we have used standard tail bounds on the binomial distribution in the Poisson regime stated in Lemma 4.6. Note that setting $x = 13$ (any $x > 12$ works) and using that $p_N \leq N^{-3/5}$ we find that if $K < 2Np_N$, then

$$\mathbb{P}(\max_i \{A_i, B_i\} > 13) \leq o(N^{-2}). \quad (4.6)$$

Assume here and below that N is large enough so that $1 - 2p_N > 2/3$. At this point we only need to show that for $\{V_1, \dots, V_K\}$ and $\{W_1, \dots, W_K\}$ two independent (unordered) collections of i.i.d. $\text{Bi}(N - 2K - 1, p_N)$ random variables with $Np_N/3 < Np_N(1 - 2p_N)/2 < K < 2Np_N$, and for fixed non-negative $A_1, \dots, A_K, B_1, \dots, B_K$ such that each A_i and B_i are no greater than 13,

$$\mathbb{P}(\{V_1 + A_1, \dots, V_K + A_K\} = \{W_1 + B_1, \dots, W_K + B_K\}) = o(1/N^2). \quad (4.7)$$

Rather than dealing with the multi-sets, we look instead at the (nearly multinomial) vectors of counts. For $i = -13, \dots, N - 2K + 12$, let $X_i = |\{j : V_j + A_j = i\}|$ be the number of the $(V_j + A_j)$'s that are equal to i and $Y_i = |\{j : W_j + B_j = i\}|$ be the analogous counts for the $(W_j + B_j)$'s. The left hand side of (4.7) is equal to

$$\mathbb{P}(X_j = Y_j, j = -13, \dots, N - 2K + 12) \leq \mathbb{P}(X_{j_i} = Y_{j_i}, i = 0, \dots, \lfloor \alpha\sqrt{Np_N} \rfloor - 1), \quad (4.8)$$

where $\alpha > 0$ will be chosen later and we define $j_i = \lfloor Np_N \rfloor + i$. To shorten formulas define the index set $\mathcal{I} = \mathcal{I}(N, \alpha) := \{0, \dots, \lfloor \alpha\sqrt{Np_N} \rfloor - 1\}$. We bound the probability (4.8) by showing first that for an appropriate $\delta > 0$,

$$\mathbb{P}(X_{j_i} > (1 + \delta)\mathbb{E}X_{j_i}, \text{ for some } i \in \mathcal{I}) = o(N^{-2}), \quad (4.9)$$

and then that given $X_{j_i} \leq (1 + \delta)\mathbb{E}X_{j_i}$ for all $i \in \mathcal{I}$, we apply the local central limit theorem to the Y_{j_i} (represented as sums of independent Bernoulli variables) to show that the event on the right hand side of (4.8) has chance $o(N^{-2})$.

To show (4.9), first note that by the local central limit theorem for the binomial distribution (noting that $p_N \rightarrow 0$), there are positive constants $c_1 = c_1(\alpha)$ and c_2 such that for all $i \in \mathcal{I}$ and $k = 1, \dots, K$,

$$\frac{c_1}{\sqrt{Np_N}} \leq \mathbb{P}(V_k + A_k = j_i), \mathbb{P}(W_k + B_k = j_i) \leq \frac{c_2}{\sqrt{Np_N}}. \quad (4.10)$$

Thus, for each j_i , X_{j_i} is a sum of K independent Bernoulli variables, each having success probability upper and lower bounded as per (4.10), and using this, a union bound, Lemma 4.6, and the bounds on K and p_N , we have

$$\begin{aligned} \mathbb{P}(X_{j_i} > (1 + \delta)\mathbb{E}X_{j_i}, \text{ for some } i \in \mathcal{I}) \\ &\leq \sum_{i \in \mathcal{I}} 2 \exp \left\{ -\frac{\delta^2}{2 + \delta} \mathbb{E}X_{j_i} \right\} \\ &\leq 2\alpha \sqrt{Np_N} \exp \left\{ -\frac{\delta^2}{2 + \delta} \frac{c_1}{3} \sqrt{Np_N} \right\} \\ &\leq 2\alpha N^{1/5} N^{-\frac{\delta^2}{2+\delta} \frac{c_1}{3}} o(1). \end{aligned}$$

Now choosing

$$\delta = \frac{1 + \sqrt{1 + 40c_1/27}}{10c_1/27},$$

shows (4.9) is satisfied, since for this choice of δ ,

$$-\frac{\delta^2}{2 + \delta} \frac{c_1}{3} + 1/5 = -2.$$

To finish the proof, we show that for an appropriate choice of α (small),

$$\mathbb{P}(X_{j_i} = Y_{j_i}, i \in \mathcal{I} | X_{j_i} \leq (1 + \delta)\mathbb{E}X_{j_i}, \text{ all } i \in \mathcal{I}) = o(N^{-2}).$$

Let $K_0 = K$ and $K_i = K - \sum_{\ell=0}^{i-1} Y_{j_\ell}$ and define \mathcal{F}_i to be the sigma field generated by Y_{j_0}, \dots, Y_{j_i} . Observe that for each $i \in \mathcal{I}$, given \mathcal{F}_{i-1} , Y_{j_i} is a sum of K_i Bernoulli variables, each having success probability Q satisfying (using (4.10))

$$\frac{c_1}{\sqrt{Np_N}} \leq \frac{c_1/\sqrt{Np_N}}{1 - ic_1/\sqrt{Np_N}} \leq Q \leq \frac{c_2/\sqrt{Np_N}}{1 - ic_2/\sqrt{Np_N}} \leq \frac{c_2}{\sqrt{Np_N}} (1 - \alpha c_2)^{-1}.$$

So we demand that $(1 - \alpha c_2) > 0$ which is not an issue: changing α affects only c_1 and δ in the argument above. Moreover, by decreasing α , we increase c_1 , and as $\alpha \rightarrow 0$, c_1 stays bounded from above (since it's no greater than c_2) and thus so does δ . The local central limit for sums of independent Bernoulli variables implies that

$$\begin{aligned} \mathbb{P}(Y_{j_i} = X_{j_i} | X_{j_i} \leq (1 + \delta)\mathbb{E}X_{j_i} \text{ all } i \in \mathcal{I}; Y_{j_\ell} = X_{j_\ell} \text{ all } \ell = 0, \dots, i-1; \mathcal{F}_{i-1}) \\ \leq C \left[K_i \frac{c_1}{\sqrt{Np_N}} \left(1 - \frac{c_2}{\sqrt{Np_N}} (1 - \alpha c_2)^{-1} \right) \right]^{-1/2}, \end{aligned} \quad (4.11)$$

for some constant C . Now the condition that $X_{j_i} \leq (1 + \delta)\mathbb{E}X_{j_i}$ and the lower bound on K implies that

$$\begin{aligned} K_i &\geq \frac{Np_N}{3} - (1 + \delta) \sum_{\ell=0}^{i-1} \mathbb{E}X_{j_\ell} \\ &\geq \frac{Np_N}{3} - (1 + \delta) \sum_{\ell \in \mathcal{I}} \mathbb{E}X_{j_\ell} \geq \frac{Np_N}{3} - (1 + \delta)\alpha 2c_2 Np_N, \end{aligned}$$

where we have used that $\mathbb{E}X_{j_i} \leq Kc_2/\sqrt{Np_N} \leq 2c_2\sqrt{Np_N}$. By choosing α small enough (so that $1/3 - 2c_2(1 + \delta)\alpha > 0$) we find that K_i is at least of order Np_N for all $i \in \mathcal{I}$ so that (4.11) is $O((Np_N)^{-1/4})$. Now moving through \mathcal{I} sequentially, we have

$$\begin{aligned} \mathbb{P}(X_{j_i} = Y_{j_i}, i \in \mathcal{I} | X_{j_i} \leq (1 + \delta)\mathbb{E}X_{j_i}, \text{ all } i \in \mathcal{I}) \\ = \exp \left\{ -\frac{\alpha}{4}\sqrt{Np_N}(\log(Np_N) + O(1)) \right\} \\ \leq \exp \left\{ -\frac{\alpha}{4}\log(N)(\log(\log(N)) + O(1)) \right\} = o(N^{-2}). \end{aligned}$$

□

Lemma 4.6. *Let X be the sum of independent indicators. Then for any $\varepsilon > 0$,*

$$\begin{aligned} \mathbb{P}(X \leq \mathbb{E}X(1 - \varepsilon)) &\leq \exp \left\{ -\frac{\varepsilon^2}{2}\mathbb{E}X \right\}, \\ \mathbb{P}(X \geq \mathbb{E}X(1 + \varepsilon)) &\leq \exp \left\{ -\frac{\varepsilon^2}{2 + \varepsilon}\mathbb{E}X \right\}. \end{aligned}$$

If X is a binomial distribution and $x > 0$, then

$$\mathbb{P}(X > x) \leq \left(\frac{e}{x}\right)^x (\mathbb{E}X)^x.$$

Proof. The first statement is a standard Chernoff bound for sums of independent indicators. The second follows in the usual way but we prove this particular form. For any $\theta > 0$, a direct computation yields

$$\mathbb{P}(X > x) \leq e^{-\theta x} \mathbb{E}e^{\theta X} \leq \exp \left\{ \mathbb{E}X(e^\theta - 1) - \theta x \right\}.$$

Setting $\theta = \log(1 + x/\mathbb{E}X)$ in the previous formula and simplifying yields

$$\mathbb{P}(X > x) \leq \left(\frac{e}{x + \mathbb{E}X}\right)^x (\mathbb{E}X)^x \leq \left(\frac{e}{x}\right)^x (\mathbb{E}X)^x,$$

as desired. □

We finish the section with a couple open problems. In Theorem 4.5 is it possible to identify from 2-neighborhoods? What happens in the regime of p_N we don't handle, where $\omega(N^{-1}) = p_N = O(\log(N)^2/N)$?

5 The Random Jigsaw Puzzle

Consider a factory that manufactures jigsaw puzzles - with the goal of producing individual unique puzzles that can be assembled. Since the images on the puzzle might not be informative (e.g. if there is a large patch of sky), the factory aims to make sure that a unique assembly of the puzzle is guaranteed just from the shape of the interface of the pieces. Assume that there are q different type of interfaces which we call “jigs” and the puzzle is of size $n \times n$. How large should q be so that the puzzle can be uniquely assembled? Note that intuitively assembly of the puzzle is harder the smaller q is. In this section we provide upper and lower bounds on q in terms of n to determine identifiability. The scaling between q and n is stated as an open problem.

5.1 Formal description of the model

We use a more intuitive description than that of Example 3 in the introduction. The puzzle is given by an $n \times n$ grid of squares where adjacent squares share an edge. Each *edge* of a square in the grid is colored uniformly at random from one of q colors. A *piece* of the puzzle consists of a “vertex” at the center of the square along with the four adjacent colored edges. Vertices at the edge of the grid have an edge on the border of the grid so that each vertex has exactly four edges associated to it. Given two pieces both having an edge of the same color, we assume that there is a unique way to connect the two pieces (i.e. there are no symmetries in the jigs). The input to the problem is all pieces and the desired output is the original composition of the puzzle.

We first use blocking configurations to obtain an easy negative result.

Proposition 5.1. *If $q = o(n^{2/3})$ then the probability of identification goes to 0 as $n \rightarrow \infty$.*

Proof. Call a pair of piece *aligned* if it is at position $(j, 2i), (j, 2i + 1)$. Let $X_{i,j,i',j'}$ be the indicator of the following event. Consider the map $\pi : (x, y) \rightarrow (x - j + j', y - 2i + 2i')$. Let $X_{i,j,i',j'}$ be 1 if all edges emanating from $(j, 2i), (j, 2i + 1)$ have the same color as their π images except that the edge connecting $(j, 2i)$ and $(j, 2i + 1)$ has a different color than its image under π . Note that if $X_{i,j,i',j'} = 1$ then there isn't a unique solution to the puzzle as the two aligned parts can be exchanged. Note that here use the fact that with high probability there are no automorphism of the labelled puzzle (even excluding two neighborhoods). Let

$$Y = \sum_{(i,j) \neq (i',j')} X_{i,j,i',j'}.$$

Then $\mathbb{E}X_{i,j,i',j'} = q^{-6}(1 - 1/q)$ and moreover, it is easy to check that the $X_{i,j,i',j'}$ are pairwise independent. Thus

$$\text{Var}(Y) = \sum_{(i,j) \neq (i',j')} \text{Var}(X_{i,j,i',j'}) \leq \mathbb{E}[Y]$$

It follows that if $n^4 q^{-6} \rightarrow \infty$ then $\mathbb{E}[Y] \rightarrow \infty$ and so by the second moment method, $\mathbb{P}[Y \geq 1] \rightarrow 1$, concluding the proof. \square

On the other hand, if $q \gg n^4$, then by considering expectations, the number of edges with the same color tends to zero in probability and identification is trivial, so if $q = \omega(n^4)$, then the probability of identification tends to 1 as $n \rightarrow \infty$. In fact we can do better.

Proposition 5.2. *If $q = \omega(n^2)$ then it is possible to assemble the puzzle with probability tending to one. More formally, if $q = \omega(n^2)$ then there exists an algorithm such that the probability it correctly assembles the puzzle (up to rotations) tends to one.*

Proof. We show that with probability tending to one, we can assemble the puzzle by first joining edges with colors that appear exactly once in the puzzle and then filling in any remaining holes. Write $q = 2cn(n+1)$ and let $m = 2n(n+1)$ be the number of edges. Let U be the number of colors which appear exactly once. Then

$$\mathbb{E}U = q \frac{m}{q} (1 - 1/q)^{m-1} \geq m(1 - 1/c).$$

Also note that U is a function of the independent edge colors such that if a single color changes, then U can change by at most 2. Thus we can apply McDiarmid's inequality for bounded differences to obtain that

$$\mathbb{P}(U \geq m(1 - 2/c)) \geq 1 - \exp \left\{ \frac{-m}{2c^2} \right\}.$$

Given U , the locations of the edges that receive unique colors is exchangeable and so on $U \geq m(1 - 2/c)$, U dominates the Bernoulli- $(1 - 3/c)$ product measure on edges with chance at least $1 - \exp\{-m(1 - 2/c)^2/(3 - 1/c)\}$, using, e.g., the Chernoff bound of Lemma 4.6. Thus, on the good event that $U \geq m(1 - 2/c)$ and at most $m(1 - 2/c)$ of the Bernoulli variables are 1, we can generate the locations of the unique colored edges by first generating the Bernoulli variables on edges and then adding the appropriate number of unique colors to the remaining edges chosen uniformly at random.

If c is large enough so that $1 - 3/c > 0.9$ (say), then standard results in percolation theory [Grimmett, 1999, (8.97-8)] imply that the graph induced by the positive Bernoulli variables in the box (which on the good event are dominated by the unique edge color indicators) has a connected component touching all boundaries. Once such a component is determined (up to rotations), it is not hard to complete the puzzle. By considering expectations, the probability of having two pieces that share two or more colors tends to zero. Thus given a location of a piece neighboring two pieces that are already assembled – i.e., an empty corner – there is a unique piece that can fit there.

Consider the process of starting with component formed by joining edges with unique colors and then repeatedly adding pieces to vacant corners. With probability tending to one, when this process terminates, the collection of vertices covered has no empty corners. It is easy to see that this implies that the complete puzzle has been recovered. \square

Remark 5.3. We have assumed that “edge” pieces of the puzzle cannot be distinguished from interior pieces. If the edge pieces can be distinguished, then the proposition still holds since with probability tending to one it is possible to construct the border by matching colors that only appear once on the border and then filling in the interior using corners as is done in the proof above. It’s interesting that without the border, we need a non-trivial result from percolation theory to start the algorithm.

6 Conclusion and Additional Open Problems

A number of open problems regarding sharper bounds and extension to other models are mentioned in the text and can be summarized as follows:

Problem 6.1. *For the graph shotgun problem on boxes in \mathbb{Z}_n^d with labels given by i.i.d., Ising, Potts model, proper coloring etc., find the threshold for the graph identification problem.*

It is natural to consider canonical fixed graphs other than the lattice. As illustrated in the introduction, the case of regular trees should be rather straightforward for many of these models. However, other families of graphs may be amenable to analysis, e.g., expander graphs.

Problem 6.2. *For the graph shotgun problem on a random graph model, e.g., Erdős-Rényi, preferential attachment, configuration, random regular graphs, etc., find the threshold for the graph identification problem.*

This question applies to both the labeled and unlabeled case. It is also interesting to understand if the graph identification problem shares properties of other constraint satisfaction problems:

Problem 6.3. *Are there graph shotgun problems for which there is a “computationally hard” but identifiable regime.*

This problem identifies graph shotgun assembly as a constraint satisfaction problem: for each neighborhood we have to find all intersecting neighborhoods. In the language of constraint satisfaction, the problem would be classified as *planted*, meaning that we start from a solution and then impose constraints based on the solution.

Acknowledgments

E.M. would like to acknowledge the support of the following grants: NSF grants DMS 1106999 and CCF 1320105, DOD ONR grant N00014-14-1-0823, and grant 328025 from the Simons Foundation. N.R. received support from ARC grant DP150101459 and thanks Aslan Tchamkerten for helpful discussions about DNA shotgun assembly. We thank James Lee for suggesting the terminology “jigs” and a reviewer for pointers to the graph isomorphism problem literature.

References

- [Addario-Berry et al., 2012] Addario-Berry, L., Broutin, N., and Goldschmidt, C. (2012). The continuum limit of critical random graphs. *Probab. Theory Related Fields*, 152(3-4):367–406.
- [Aldous and Lyons, 2007] Aldous, D. and Lyons, R. (2007). Processes on unimodular random networks. *Electron. J. Probab.*, 12:no. 54, 1454–1508.
- [Arratia et al., 1996] Arratia, R., Martin, D., Reinert, G., and Waterman, M. S. (1996). Poisson process approximation for sequence repeats, and sequencing by hybridization. *J. Comp. Bio.*, 3(3):425–463.
- [Babai et al., 1980] Babai, L., Erdős, P., and Selkow, S. M. (1980). Random graph isomorphism. *SIAM J. Comput.*, 9(3):628–635.
- [Bollobás, 1981] Bollobás, B. (1981). The diameter of random graphs. *Trans. Amer. Math. Soc.*, 267(1):41–52.
- [Cai et al., 1992] Cai, J.-Y., Fürer, M., and Immerman, N. (1992). An optimal lower bound on the number of variables for graph identification. *Combinatorica*, 12(4):389–410.
- [Dyer et al., 1994] Dyer, M., Frieze, A., and Suen, S. (1994). The probability of unique solutions of sequencing by hybridization. *J. Comp. Bio.*, 1(2):105–110.
- [Frisch and Tamuz, 2014] Frisch, J. and Tamuz, O. (2014). Transitive graphs uniquely determined by their local structure. Preprint <http://arxiv.org/abs/1411.6534>.
- [Grimmett, 1999] Grimmett, G. (1999). *Percolation*, volume 321 of *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*. Springer-Verlag, Berlin, second edition.
- [Harary, 1974] Harary, F. (1974). A survey of the reconstruction conjecture. In *Graphs and combinatorics (Proc. Capital Conf., George Washington Univ., Washington, D.C., 1973)*, pages 18–28. Lecture Notes in Math., Vol. 406. Springer, Berlin.
- [Kelly, 1957] Kelly, P. J. (1957). A congruence theorem for trees. *Pacific J. Math.*, 7:961–968.
- [Łuczak, 1998] Łuczak, T. (1998). Random trees and random graphs. In *Proceedings of the Eighth International Conference “Random Structures and Algorithms” (Poznan, 1997)*, volume 13, pages 485–500.
- [Motahari et al., 2013] Motahari, A. S., Bresler, G., and Tse, D. N. (2013). Information theory of DNA shotgun sequencing. *Information Theory, IEEE Transactions on*, 59(10):6273–6289.

- [Nachmias and Peres, 2008] Nachmias, A. and Peres, Y. (2008). Critical random graphs: diameter and mixing time. *Ann. Probab.*, 36(4):1267–1286.
- [Pitman, 1997] Pitman, J. (1997). Probabilistic bounds on the coefficients of polynomials with only real zeros. *J. Combin. Theory Ser. A*, 77(2):279–303.
- [Riordan and Wormald, 2010] Riordan, O. and Wormald, N. (2010). The diameter of sparse random graphs. *Combin. Probab. Comput.*, 19(5-6):835–926.
- [Soudry et al., 2013] Soudry, D., Keshri, S., Stinson, P., Oh, M.-h., Iyengar, G., and Paninski, L. (2013). A shotgun sampling solution for the common input problem in neural connectivity inference. Preprint <http://arxiv.org/abs/1309.3724>.